

# Matrix Completion under Interval Uncertainty

Jakub Mareček, Peter Richtárik, Martin Takáč \*

April 4, 2016

## Abstract

Matrix completion under interval uncertainty can be cast as matrix completion with element-wise box constraints. We present an efficient alternating-direction parallel coordinate-descent method for the problem. We show that the method outperforms any other known method on a benchmark in image in-painting in terms of signal-to-noise ratio, and that it provides high-quality solutions for an instance of collaborative filtering with 100,198,805 recommendations within 5 minutes.

## 1 Motivation

Matrix completion is a well-known problem, with applications ranging from image processing to recommender systems. When dimensions of a matrix  $X$  and some of its elements  $X_{i,j}, (i,j) \in \mathcal{E}$  are known, the goal is to find the unknown elements. Without imposing any further requirements on  $X$ , there are infinitely many solutions. In many applications, however, the matrix completion that minimizes the rank:

$$\min_Y \text{rank}(Y) \quad \text{subject to} \quad Y_{i,j} = X_{i,j}, \quad (i,j) \in \mathcal{E}, \quad (1)$$

provides the simplest explanation for the data. There is a long history of work on the problem, c.f. [9, 36, 47, 21], with thousands of papers published annually since 2010. We hence cannot provide a complete overview.

Let us note that Fazel [10] suggested to replace the rank, which is the sum of non-zero elements of the spectrum, with the nuclear norm, which is the sum of the spectrum. The minimisation of the nuclear norm can be cast as a semidefinite programming (SDP) problem and approaches based on the nuclear-norm have proven very successful in theory [6] and very popular in practice. [36, 3] study the Singular Value Thresholding (SVT) algorithm. This, however, required the computation of a singular value decomposition (SVD) in each iteration. A number of other approaches, e.g., augmented Lagrangian methods [44],

---

\*Jakub Mareček is at IBM Research, Martin Takáč is at Lehigh University, and Peter Richtárik is at the University of Edinburgh. Their addresses are [jakub@marecek.cz](mailto:jakub@marecek.cz), [takac.mt@gmail.com](mailto:takac.mt@gmail.com), and [peter.richtarik@ed.ac.uk](mailto:peter.richtarik@ed.ac.uk), respectively.

appeared, but those would require a truncated SVD or a number of iterations [15, 22, 37, 45] of the power method. Even considering the recent progress in randomized methods for approximating SVD, [13], the approximation becomes very time-consuming as the dimensions of matrices grow.

A major computational break-through came in the form of the alternating least squares (ALS) algorithms [40, 31]. Initially, the algorithm has been used as a heuristic for finding stationary points of the non-convex problem [40, 31, 27, 2, 12], where a single iteration had complexity  $O(|\mathcal{E}|r^2)$ , for  $|\mathcal{E}|$  observations and rank  $r$ , c.f., p. 60 in [20]. Keshavan et al. [19, 20], however, proved its exponential rate of convergence to the global optimum with high probability, under probabilistic assumptions common in the compressed sensing community. Further, more technical analyses of the convergence to the global optimum have been performed by Jain et al. [17].

Many studies of matrix completion consider the uncertainty, in some form. A number of analyses [19, 20, 17] consider the use of the standard rank-minimisation for the reconstruction of low-rank  $m \times n$  matrix  $XY^T$  from  $XY^T + W$ , where  $X \in \mathbb{R}^{m \times r}$ ,  $Y \in \mathbb{R}^{n \times r}$ ,  $W \in \mathbb{R}^{m \times n}$  with elements of  $W$  being bounded i.i.d. random variables, which are sub-Gaussian and have bounded expectation. A number of further analyses [46, 4] considered the use of the standard rank-minimisation for the reconstruction of low-rank  $m \times n$  matrix  $XY^T$  from  $XY^T + S$ , where  $X, Y$  are as above and  $W$  has a small number of non-zero entries. [7] consider some columns being corrupted. Although we are not aware of any studies of matrix completion under interval uncertainty, interval-based uncertainty has been considered in related problems. Alaiz et al. [1] consider the min-max variant of the problem of finding the nearest correlation matrix, i.e., the problem of finding the closest matrix within the set of symmetric positive definite matrices with the unit diagonal to an uncertainty set, with respect to the Frobenius norm. [23] studied interval uncertainty in certain semidefinite programming problems, which can be used to encode the nuclear-norm minimisation.

In contrast, we present an extension of matrix completion toward interval uncertainty, which has applications in image inpainting, collaborative filtering, and beyond. The algorithm we present for solving the problem can be seen as a coordinate-wise version of the ALS algorithm, which does not require the approximation of the spectrum of the matrix. This makes it possible solve complete matrices  $480,189 \times 17,770$  matrix within minutes on a standard laptop. First, we provide an overview of the possible applications.

## 1.1 Collaborative Filtering under Uncertainty

Collaborative filtering is a well-established application of matrix completion problems [39], largely thanks to the success of the Netflix Prize. There is a matrix, where each row corresponds to one user and each column corresponds to a product or service. Considering that every user rates only a modest number of products or services, there are only a small number of entries of the matrix known. Our extension is motivated by the fact, that one user may provide two

different ratings for one and the same product at two different times, depending on the current mood and other circumstances at the two times. One may hence want to consider an interval  $[\underline{x}, \bar{x}]$  instead of a fixed value  $x$  of the rating, e.g.,  $[x - \epsilon, x + \epsilon]$ . Further, when one knows the scale  $[0, M]$  the rating  $x$  is chosen from, one can consider  $[\max\{0, x - \epsilon\}, \min\{x + \epsilon, M\}]$ . Hence, if intervals are known for elements  $X_{i,j}$  of a matrix  $X$  indexed by  $(i, j) \in \mathcal{I}$ , one may want to solve:

$$\begin{aligned} \min_{Y_{i,j} \in [0, M]} \max_{X_{i,j} \in [\underline{X}_{i,j}, \overline{X}_{i,j}] \forall (i,j) \in \mathcal{I}} \text{rank}(Y) \\ \text{subject to } Y_{i,j} = X_{i,j}, \quad \forall (i, j) \in \mathcal{I}. \end{aligned} \quad (2)$$

Although numerous extensions of matrix completion problems have been studied, e.g. [26], the use of robustness to interval uncertainty is novel. It can be seen as an extension of robust optimisation [38] to matrix completion.

## 1.2 Image In-Painting

Further applications can be found in image processing. In in-painting problems, a subset of pixels from an image are given and the goal is to fill in the missing pixels. Rank-constrained matrix completion with equalities, where  $\mathcal{I}$  is the index set of all known pixels, has been used numerous times [6, 16, 25, 11, 22, 15, 45, 45] in this setting. If the image comes from real sensors, it the corresponding matrix may have full (numerical) rank, but have quickly decreasing singular values in its spectrum. In such a case, instead of solving the equality-constrained problem (1), one should like to find a low-rank approximation  $Y^*$  of  $X$ , such that the known entry of  $X$  is not far away from  $Y^*$ , i.e.,  $\forall (i, j) \in \mathcal{I}$  we have  $Y_{i,j} \approx X_{i,j}$ . Let us illustrate this with a small matrix

$$X = \begin{pmatrix} 68.16 & 78.12 & 24.04 \\ 78.12 & 90.09 & 30.03 \\ 24.04 & 30.03 & 20.01 \end{pmatrix},$$

which has rank 3 and its singular values  $\Sigma = (167.9945, 10.2553, 0.0102)^T$ . It is easy to verify that

$$Y^*(2) = \begin{pmatrix} 68.1546 & 78.1250 & 24.0389 \\ 78.1250 & 90.0853 & 30.0310 \\ 24.0389 & 30.0310 & 20.0098 \end{pmatrix}$$

is the best rank 2 approximation of  $X$  in Frobenius norm. Observe that no single element of  $Y^*(2)$  is identical to  $X$ , but that  $Y^*(2) \approx X$ . It is an easy exercise to show that for any  $X \in \mathbb{R}^{m \times n}$  with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m,n\}}$ , and  $Y^*(r)$  as its best rank- $r$  approximation, we have  $|X_{i,j} - (Y^*(r))_{i,j}| \leq \sum_{i=r+1}^{\min\{m,n\}} \sigma_i =: \mathcal{R}(r)$  for all  $(i, j)$ . Therefore, one should not require equality constraints in (1), but rather inequalities  $|Y_{i,j} - X_{i,j}| \leq \mathcal{R}(r), \forall (i, j) \in \mathcal{I}$ . Notice that this approach is not the same as minimizing  $\sum_{(i,j) \in \mathcal{I}} (X_{i,j} - Y_{i,j})^2$  over all rank  $r$  matrices, because we do not penalize the elements of  $Y$ , which

are already close to  $X$ . It is also different from the usual treatment of noise in the observations [5]. One could rather formulate this as the minimization of  $\sum_{(i,j) \in \mathcal{I}} \max\{0, |X_{i,j} - Y_{i,j}| - \mathcal{R}(r)\}^2$  over all rank  $r$  matrices. Further, one knows the range of values allowed, e.g.,  $[0, 1]$  for common encoding of gray-scale images. This can hence be seen as “side information” which, as we will show in numerical section, improves recovery of a low-rank approximation considerably. Further still, one could assume that the intensity should be at least 0.8, if pixels are missing within a light region of the image, or similar domain-specific heuristics.

A number of other applications, e.g., in the recovery of structured matrices [8], forecasting with side information, and in sparse principal component analysis with priors on the principal components, can be envisioned.

## 2 The Problem

Formally, let  $X$  be an  $m \times n$  matrix to be reconstructed. Assume that elements  $(i, j) \in \mathcal{E}$  of  $X$  we wish to fix, for elements  $(i, j) \in \mathcal{L}$  we have lower bounds and for elements  $(i, j) \in \mathcal{U}$  we have upper bounds. We employ the following natural formulation for the equality and inequality constrained matrix completion problem:

$$\begin{aligned} \min_{X \in \mathbb{R}^{m \times n}} \quad & \text{rank}(X) \\ \text{subject to} \quad & X_{ij} = X_{ij}^{\mathcal{E}}, \quad (i, j) \in \mathcal{E} \\ & X_{ij} \geq X_{ij}^{\mathcal{L}}, \quad (i, j) \in \mathcal{L} \\ & X_{ij} \leq X_{ij}^{\mathcal{U}}, \quad (i, j) \in \mathcal{U}. \end{aligned} \tag{3}$$

We shall enforce the following natural assumption:

**Assumption 1.**  $\mathcal{E} \cap (\mathcal{L} \cup \mathcal{U}) = \emptyset$  and  $X_{ij}^{\mathcal{L}} \leq X_{ij}^{\mathcal{U}}$  whenever  $(ij) \in \mathcal{L} \cap \mathcal{U}$ .

The first condition says that if some element  $(ij)$  is already fixed by an equality constraint, it does not (unnecessarily) appear any of the inequality constraints. The second condition says the upper and lower bounds should be consistent.

Problem (3) is NP-hard, even with  $\mathcal{U} = \mathcal{L} = \emptyset$  [28, 14]. A number of special cases of (3) have been studied in the literature, e.g., in [36, 30, 18]. A popular heuristic enforces low rank in a synthetic way by writing  $X$  as a product of two matrices,  $X = LR$ , where  $L \in \mathbb{R}^{m \times r}$  and  $R \in \mathbb{R}^{r \times n}$ . Hence,  $X$  is of rank at most  $r$  [41]. Let  $L_{i\cdot}$  and  $R_{\cdot j}$  be the  $i$ -th row and  $j$ -th column of  $L$  and  $R$ , respectively. Instead of (3), we consider the *smooth, non-convex* problem

$$\min\{f(L, R) : L \in \mathbb{R}^{m \times r}, R \in \mathbb{R}^{r \times n}\}, \tag{4}$$

where

$$f(L, R) := \frac{\mu}{2} \|L\|_F^2 + \frac{\mu}{2} \|R\|_F^2 + f_{\mathcal{E}}(L, R) + f_{\mathcal{L}}(L, R) + f_{\mathcal{U}}(L, R). \tag{5}$$

---

**Algorithm 1** MACO: Matrix Completion via Alternating Parallel Coordinate Descent

---

Input:  $\mathcal{E}, \mathcal{L}, \mathcal{U}, X^{\mathcal{E}}, X^{\mathcal{L}}, X^{\mathcal{U}}$ , rank  $r$   
 Output:  $m \times n$  matrix  $LR$   
 1: choose  $L \in \mathbb{R}^{m \times r}$  and  $R \in \mathbb{R}^{r \times n}$   
 2: **for**  $k = 0, 1, 2, \dots$  **do**  
 3:   choose random subset  $\hat{\mathcal{S}}_{\text{row}} \subset \{1, \dots, m\}$   
 4:   **for**  $i \in \hat{\mathcal{S}}_{\text{row}}$  **in parallel do**  
 5:     choose  $\hat{r} \in \{1, \dots, r\}$  uniformly at random  
 6:     compute  $\delta_{i\hat{r}}$  using formula (8)  
 7:     update  $L_{i\hat{r}} \leftarrow L_{i\hat{r}} + \delta_{i\hat{r}}$   
 8:   **end for**  
 9:   choose random subset  $\hat{\mathcal{S}}_{\text{column}} \subset \{1, \dots, n\}$   
 10:   **for**  $j \in \hat{\mathcal{S}}_{\text{column}}$  **in parallel do**  
 11:     choose  $\hat{r} \in \{1, \dots, r\}$  uniformly at random  
 12:     compute  $\delta_{\hat{r}j}$  using (11)  
 13:     update  $R_{\hat{r}j} \leftarrow R_{\hat{r}j} + \delta_{\hat{r}j}$   
 14:   **end for**  
 15: **end for**

---

Above we have

$$\begin{aligned}
 f_{\mathcal{E}}(L, R) &:= \frac{1}{2} \sum_{(ij) \in \mathcal{E}} (L_{i:} R_{:j} - X_{ij}^{\mathcal{E}})^2 \\
 f_{\mathcal{L}}(L, R) &:= \frac{1}{2} \sum_{(ij) \in \mathcal{L}} (X_{ij}^{\mathcal{L}} - L_{i:} R_{:j})_+^2 \\
 f_{\mathcal{U}}(L, R) &:= \frac{1}{2} \sum_{(ij) \in \mathcal{U}} (L_{i:} R_{:j} - X_{ij}^{\mathcal{U}})_+^2,
 \end{aligned}$$

where  $\xi_+ = \max\{0, \xi\}$ .

The parameter  $\mu > 0$  helps to prevent scaling issues<sup>1</sup>. We could optionally set  $\mu$  to zero and instead, from time to time, rescale matrices  $L$  and  $R$ , so that their product is not changed [41]. The term  $f_{\mathcal{E}}$  (resp.  $f_{\mathcal{U}}$ ,  $f_{\mathcal{L}}$ ) encourages the equality (resp. inequality) constraints to hold.

### 3 The Method

Coordinate descent algorithms (CDA) are effective in solving large-scale problems, due to their low per-iteration computational cost. Although each iteration of CDA is cheap, many more iterations are required for convergence, compared to second-order algorithms or similar. Recently, the stochastic CDA has received much attention [29, 32] not least due to the parallelizability [35, 33, 42, 43] with almost linear speed-up in regimes with sparse data, when the number of parallel updates  $\tau$  is much smaller than the dimension of the optimization problem [30]. Distributed variants have also been studied [24, 34].

---

<sup>1</sup>Let  $X = LR$ , then also  $X = (cL)(\frac{1}{c}R)$  as well, but we see that for  $c \rightarrow 0$  or  $c \rightarrow \infty$  we have  $\|L\|_F^2 + \|R\|_F^2 \ll \|cL\|_F^2 + \|\frac{1}{c}R\|_F^2$ .

In Algorithm 1, we present our alternating parallel coordinate descent method for MATRIX COMPLETION, henceforth simply “MACO”. In Steps 3–8 of our algorithm, we fix  $R$ , choose random  $\hat{r}$  and a random set  $\hat{\mathcal{S}}_{\text{row}}$  of rows of  $L$ , and update, in parallel for  $i \in \hat{\mathcal{S}}_{\text{row}}$ :  $L_{i\hat{r}} \leftarrow L_{i\hat{r}} + \delta_{i\hat{r}}$ . In Steps 9–14, we fix  $L$ , choose random  $\hat{r}$  and a random set  $\hat{\mathcal{S}}_{\text{column}}$  of columns of  $R$ , and update, in parallel for  $j \in \hat{\mathcal{S}}_{\text{column}}$ :  $R_{\hat{r}j} \leftarrow R_{\hat{r}j} + \delta_{\hat{r}j}$ .

Let us now comment on the computation of the updates,  $\delta_{i\hat{r}}$  and  $\delta_{\hat{r}j}$ . First, note that while  $f$  is not convex jointly in  $(L, R)$ , it is convex in  $L$  for fixed  $R$  and in  $R$  for fixed  $L$ .

### 3.1 Row Update

If we now fix row  $i \in \{1, 2, \dots, m\}$  and  $\hat{r} \in \{1, 2, \dots, r\}$ , and view  $f$  as a function of  $L_{i\hat{r}}$  only, it has a Lipschitz continuous derivative with constant

$$W_{i\hat{r}} = W_{i\hat{r}}(R) := \mu + \sum_{j : (ij) \in \mathcal{E}} R_{\hat{r}j}^2 + \sum_{j : (ij) \in \mathcal{L} \cup \mathcal{U}} R_{\hat{r}j}^2. \quad (6)$$

That is, for all  $L, R$  and  $\delta \in \mathbb{R}$ , we have

$$f(L + \delta E_{i\hat{r}}, R) \leq f(L, R) + \langle \nabla_L f(L, R), E_{i\hat{r}} \rangle \delta + \frac{W_{i\hat{r}}}{2} \delta^2, \quad (7)$$

where  $E_{i\hat{r}}$  is the  $n \times r$  matrix with 1 in the  $(i\hat{r})$  entry and zeros elsewhere. The minimizer of the right hand side of (7) in  $\delta$  is given by

$$\delta_{i\hat{r}} := -\langle \nabla_L f(L, R), E_{i\hat{r}} \rangle / W_{i\hat{r}}, \quad (8)$$

where

$$\begin{aligned} \langle \nabla_L f(L, R), E_{i\hat{r}} \rangle &= \mu L_{i\hat{r}} + \sum_{j : (ij) \in \mathcal{E}} (L_{i:} R_{:j} - X_{ij}^{\mathcal{E}}) R_{\hat{r}j} \\ &+ \sum_{j : (ij) \in \mathcal{U} \text{ \& } L_{i:} R_{:j} < X_{ij}^{\mathcal{U}}} (L_{i:} R_{:j} - X_{ij}^{\mathcal{U}}) R_{\hat{r}j} \\ &+ \sum_{j : (ij) \in \mathcal{L} \text{ \& } L_{i:} R_{:j} > X_{ij}^{\mathcal{L}}} (L_{i:} R_{:j} - X_{ij}^{\mathcal{L}}) R_{\hat{r}j}. \end{aligned}$$

Note that

$$f(L + \delta_{i\hat{r}} E_{i\hat{r}}, R) \leq f(L, R) - \frac{\langle \nabla_L f(L, R), E_{i\hat{r}} \rangle^2}{2W_{i\hat{r}}}. \quad (9)$$

Let  $W_{i\hat{r}}^{(k)} := W_{i\hat{r}}(R^{(k)})$  be the value of the Lipschitz constant at iteration  $k$ .

### 3.2 Column Update

Likewise, if we now fix  $\hat{r} \in \{1, 2, \dots, r\}$  and column  $j \in \{1, 2, \dots, n\}$ , and view  $f$  as a function of  $R_{\hat{r}j}$  only, it has a Lipschitz continuous derivative with constant

$$V_{\hat{r}j} = V_{\hat{r}j}(L) := \mu + \sum_{i : (ij) \in \mathcal{E}} L_{i\hat{r}}^2 + \sum_{i : (ij) \in \mathcal{U} \cup \mathcal{L}} L_{i\hat{r}}^2.$$

That is, for all  $L, R$  and  $\delta \in \mathbb{R}$ ,

$$f(L, R + \delta E_{\hat{r}j}) \leq f(L, R) + \langle \nabla_R f(L, R), E_{\hat{r}j} \rangle \delta + \frac{V_{\hat{r}j}}{2} \delta^2, \quad (10)$$

where  $E_{\hat{r}j}$  is the  $r \times m$  matrix with 1 in the  $(\hat{r}j)$  entry and zeros elsewhere. The minimizer of the right hand side of (10) in  $\delta$  is given by

$$\delta_{\hat{r}j} := -\langle \nabla_R f(L, R), E_{\hat{r}j} \rangle / V_{\hat{r}j}, \quad (11)$$

where

$$\begin{aligned} \langle \nabla_R f(L, R), E_{\hat{r}j} \rangle = & \mu R_{\hat{r}j} + \sum_{i : (ij) \in \mathcal{E}} (L_{i:R:j} - X_{ij}^{\mathcal{E}}) L_{i\hat{r}} \\ & + \sum_{i : (ij) \in \mathcal{L} \ \& \ L_{i:R:j} < X_{ij}^{\mathcal{L}}} (L_{i:R:j} - X_{ij}^{\mathcal{L}}) L_{i\hat{r}} \\ & + \sum_{i : (ij) \in \mathcal{U} \ \& \ L_{i:R:j} > X_{ij}^{\mathcal{U}}} (L_{i:R:j} - X_{ij}^{\mathcal{U}}) L_{i\hat{r}}. \end{aligned}$$

Note that

$$f(L, R + \delta_{\hat{r}j} E_{\hat{r}j}) \leq f(L, R) - \frac{\langle \nabla_R f(L, R), E_{\hat{r}j} \rangle^2}{2V_{\hat{r}j}}. \quad (12)$$

Let  $V_{\hat{r}j}^{(k)} := W_{\hat{r}j}(L^{(k)})$  be the value of the Lipschitz constant at iteration  $k$ .

### 3.3 Row and Column Sampling

The random set (“sampling”)  $\hat{\mathcal{S}}_{\text{row}}$  defined in Step 3 (resp sampling  $\hat{\mathcal{S}}_{\text{column}}$  in Step 10) can have an arbitrary distribution as long as it contains every row (resp column) of matrix  $L$  (resp  $R$ ) with positive probability. We shall now formalize this.

**Assumption 2.** *The samplings  $\hat{\mathcal{S}}_{\text{row}}$  and  $\hat{\mathcal{S}}_{\text{column}}$  are proper, i.e.,*

$$\mathbf{Prob}(i \in \hat{\mathcal{S}}_{\text{row}}) > 0 \quad \text{for all } i \in \{1, 2, \dots, m\},$$

and

$$\mathbf{Prob}(j \in \hat{\mathcal{S}}_{\text{column}}) > 0 \quad \text{for all } j \in \{1, 2, \dots, n\}.$$

In particular, we can chose the random sets  $\hat{\mathcal{S}}_{\text{row}}$  (resp  $\hat{\mathcal{S}}_{\text{column}}$ ) so that every row (resp column) has equal probability of being chosen. Samplings with this property are called *uniform*, and we use this choice in our experiments. However, our theory also allows for nonuniform samplings. If we have a multicore machine available with  $\tau$  cores, then a reasonable sampling should have cardinality  $\tau$ , or some integral multiple of  $\tau$ , so that every core has a reasonable (not too small to be underutilized, but not too large either, so as to avoid long processing time) load at every iteration.

### 3.4 The Final Step

Formulae (8) and (11) suggest that the computation of the final step is very computationally demanding. This can, however, be avoided if we define matrices  $A \in \mathbb{R}^{m \times r}$  and  $B \in \mathbb{R}^{r \times n}$  such that  $A_{iv} = W_{iv}$  and  $B_{vj} = V_{vj}$ . After each update of the solution, we can also update those matrices. Similarly, one can store sparse residuals matrices  $\Delta_{\mathcal{E}}, \Delta_{\mathcal{L}}, \Delta_{\mathcal{U}}$ , where

$$(\Delta_{\mathcal{E}})_{i,j} = \begin{cases} L_{i:R:j} - X_{ij}^{\mathcal{E}}, & \text{if } (ij) \in \mathcal{E} \\ 0, & \text{otherwise,} \end{cases}$$

and  $\Delta_{\mathcal{U}}, \Delta_{\mathcal{L}}$  are defined in similar way. Subsequently, the computation of  $\delta_{i\hat{r}}$  or  $\delta_{\hat{r}j}$  is reduced to just a few multiplications and additions.

### 3.5 Convergence Analysis

Due to the non-convex nature of (4), one has to be satisfied with convergence to a stationary point, in general.

**Theorem 1.** *Let  $\mu > 0$  and let  $(L^{(k)}, R^{(k)})$  be the (random) matrices produced by Algorithm 1 after  $k$  iterations, assuming that  $\hat{\mathbf{S}}_{\text{row}}$  and  $\hat{\mathbf{S}}_{\text{column}}$  are proper. Then for all  $k \geq 0$ ,*

$$0 \leq f(L^{(k+1)}, R^{(k+1)}) \leq f(L^{(k)}, R^{(k)}). \quad (13)$$

*That is, the method is monotonic. Moreover, with probability 1,*

$$\lim_{k \rightarrow \infty} \inf \|\nabla_L f(L^{(k)}, R^{(k)})\| = 0,$$

*and*

$$\lim_{k \rightarrow \infty} \inf \|\nabla_R f(L^{(k)}, R^{(k)})\| = 0.$$

*Proof.* From (9) and (12) we see that for all  $k$  we have

$$f(L^{(k)}, R^{(k)}) \stackrel{(9)}{\geq} f(L^{(k+1)}, R^{(k)}) \stackrel{(12)}{\geq} f(L^{(k+1)}, R^{(k+1)}) \geq 0,$$

where the last inequality follows from the fact that all parts of  $f$  defined in (5) are non-negative.

Monotonicity (13) together with the fact that  $\mu > 0$  imply that the level set

$$\Omega_0 := \{(L, R) : f(L, R) \leq f(L^{(0)}, R^{(0)})\}$$

is bounded. Now, for all  $i \in \{1, 2, \dots, m\}$ ,  $v \in \{1, 2, \dots, r\}$  and any iteration counter  $k$  we have

$$\mu \stackrel{(6)}{\leq} W_{iv}^{(k)} \stackrel{(6)}{\leq} \mu + \|R^{(k)}\|_F^2 \stackrel{(5)}{\leq} \mu + \frac{2}{\mu} f(L^{(k)}, R^{(k)}) \leq \mu + \frac{2}{\mu} f(L^{(0)}, R^{(0)}). \quad (14)$$



In the second inequality we have used Assumption 1, and in the last inequality we have used monotonicity. The same lower and upper bounds can be established for  $V_{vj}^{(k)}$ .

We shall now establish that  $\liminf \|\nabla_L f(L^{(k)}, R^{(k)})\|_F^2 = 0$  with probability 1 (the claim  $\liminf \|\nabla_R f(L^{(k)}, R^{(k)})\|_F^2 = 0$  can be proved in an analogous way). Since

$$\|\nabla_L f(L^{(k)}, R^{(k)})\|_F^2 = \sum_{i=1}^m \sum_{v=1}^r \langle \nabla_L f(L^{(k)}, R^{(k)}), E_{iv} \rangle^2,$$

it is enough to show that for  $\Delta_{iv}^{(k)} := \langle \nabla_L f(L^{(k)}, R^{(k)}), E_{iv} \rangle$  we have  $\liminf (\Delta_{iv}^{(k)})^2 = 0$  with probability 1 for all  $i \in \{1, 2, \dots, m\}$  and  $v \in \{1, 2, \dots, r\}$ . Fix any  $i$  and  $v$ . Since  $\hat{S}_{\text{row}}$  is proper, and since  $\hat{r}$  is chosen uniformly at random in each iteration, there is an infinite sequence of iterations, indexed by  $K_{iv}$ , in which the pair  $(i, v)$  is sampled.

In view of (9) and (14), for all  $k \in K_{iv}$  we have

$$f(L^{(k+1)}, R^{(k+1)}) \leq f(L^{(k+1)}, R^{(k)}) \leq f(L^{(k)}, R^{(k)}) - \frac{(\Delta_{iv}^{(k)})^2}{C},$$

where  $C = 2(\mu + \frac{2}{\mu} f(L^{(0)}, R^{(0)}))$ . Since  $f(L, U)$  is nonnegative, it must be the case that  $\sum_{k \in K_{iv}} (\Delta_{iv}^{(k)})^2$  is finite. This means that, with probability 1,  $\lim_{k \rightarrow \infty} \inf (\Delta_{iv}^{(k)})^2 = 0$ , as desired.  $\square$

## 4 Computational Results and a Discussion

We have conducted a variety of experiments. First, we present the performance in collaborative filtering, next we compare the performance in image in-painting with classical matrix completion techniques with  $\mathcal{U} \equiv \mathcal{L} \equiv \emptyset$ . We conclude with remarks on the run-time and hardware used.

### 4.1 Collaborative Filtering

In our computational testing of collaborative filtering, we start with `smallnetflix_mm`, where the training dataset contains  $c_{\text{tr}} = 3,298,163$  integers out of  $\{1, 2, 3, 4, 5\}$ , which describe how  $m = 95,526$  users rate  $n = 3,561$  movies. Second, we use a well-known data-set, which contains 100,198,805 ratings on the same scale, obtained from 480,189 users considering 17,770 products, as available from CMU<sup>2</sup>. Third, we use Yelp's Academic Dataset<sup>3</sup>, from which we have extracted a  $252,898 \times 41,958$  matrix with 1,125,458 non-zeros, again on the 1–5 scale.

<sup>2</sup> <http://www.select.cs.cmu.edu/code/graphlab/datasets/>

<sup>3</sup> [https://www.yelp.co.uk/academic\\_dataset](https://www.yelp.co.uk/academic_dataset)

Although we know some ratings exactly on `smallnetflix_mm`, we consider (4) of (3) with interval uncertainty sets of width 2:

$$\begin{aligned} Y_{i,j} &\leq \min\{5, X_{i,j} + 1\}, (i, j) \in \mathcal{I}, \\ Y_{i,j} &\geq \max\{1, X_{i,j} - 1\}, (i, j) \in \mathcal{I}. \end{aligned} \quad (15)$$

In particular, we complete a  $95526 \times 3561$  matrix of rank 2 or 3, possibly using width-2 interval uncertainty set and scale of 1 to 5 stars in the ratings. To illustrate the impact of the this change, we present the evolution of Root-Mean-Square Error (RMSE) in Figure 2 (left). Notice that an “epoch”, which is the unit on the horizontal axis, consists of  $c_{\text{tr}}$  element updates of matrix  $L$  and  $c_{\text{tr}}$  element updates of matrix  $R$ .

Let us remark that RMSE is sensitive to the choice of  $\Delta$  and the rank of the matrix we are looking for. If the underlying matrix has a higher rank than expected,  $\Delta > 0$  can lead to smaller values of RMSE. We should also note that for some fixed  $\Delta_1$  and  $\Delta_2$ , RMSE can be better with  $\Delta_1$  for a few epochs, but then get worse when compared with  $\Delta_2$ . Hence, in practice, cross validation should be used to determine suitable value of parameter  $\Delta$ .

On the Yelp data set, we have performed 10-fold cross-validation on the training set, using varying rank. As we increased the rank from 1 to 2, 4, 8, 16, 32, and 50, the average error decreased from 1.7958 to 1.8284, 1.6464, 1.4590, 1.3395, 1.2702, and 1.2454, respectively. This seems to be comparable to the best results from the 2013 Recommender Systems Challenge<sup>4</sup>, where a smaller dataset was used.

Further, one can illustrate the effects in a matrix-recovery experiment. We use random matrices  $X \in \mathbb{R}^{20 \times 20}$  of rank 8. We sample  $p\%$  of entries of the matrix and store their indices in  $\mathcal{I}$ . We solve (4) with just the inequality constraints, i.e.,  $\mathcal{E} \equiv \emptyset, \mathcal{U} \equiv \mathcal{L} \equiv \mathcal{I}$ ,  $X^{\mathcal{U}} = X - \Delta \mathbf{1}$  and  $X^{\mathcal{L}} = X + \Delta \mathbf{1}$ , where  $\mathbf{1} \in \mathbb{R}^{m \times n}$  is a matrix with all elements equals to 1. Let us denote by  $Y^*(\Delta)$  the solution of that optimization problem after  $10^5$  serial iterations ( $|\hat{\mathcal{S}}| = 1$ ) and with  $\mu = 10^{-5}$ . Figure 1 shows the dependence of error defined as follows  $\text{Error}(\Delta) = \frac{\|Y^*(\Delta) - X(7)\|_F}{\|X(7)\|_F}$ , where  $X(r)$  is the best rank  $r$  approximation of  $X$  obtain using SVD decomposition of the whole matrix. Figure 1 clearly suggest that, e.g., if 50% of elements are observed then by allowing each entry  $\in \mathcal{I}$  of reconstructed matrix to lie in  $\Delta$  neighborhood of observed values, we can decrease the relative error of reconstruction from approximately 1.22 to 0.4 for  $\Delta \approx \mathcal{R}(r)$ . In this case, the value of  $\|X(7)\|_F$  was 21.3245 and  $\mathcal{R}(r) = 0.1075$ .

## 4.2 Image In-Painting

Further, we provide a comparison on the in-painting benchmark of [45]. Table 1 details the performance of SVT [6], SVP [16], SoftImpute [25], LMaFit [11], ADMiRA, [22], JS [15], OR1MP [45], and EOR1MP [45] on 10 well-known gray-scale images (Barbara, Cameraman, Clown, Couple, Crowd, Girl, Goldhill, Lenna, Man, Peppers) of  $512 \times 512$  pixels each. 50% of pixels were removed

<sup>4</sup> <https://www.kaggle.com/c/yelp-recsys-2013>

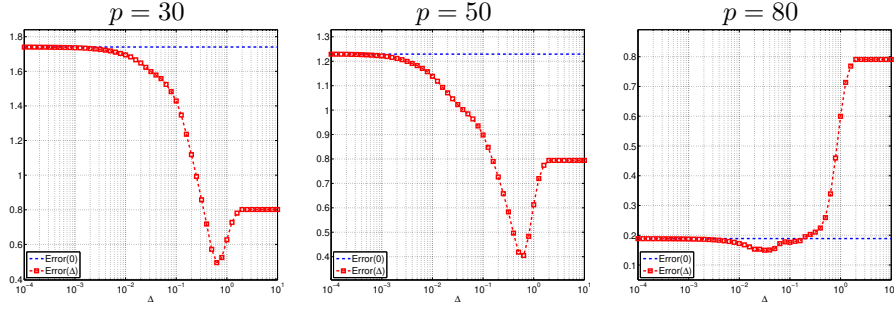


Figure 1: Dependence of Error on  $\Delta$  for various  $p \in \{30, 50, 80\}$  in matrix reconstruction.

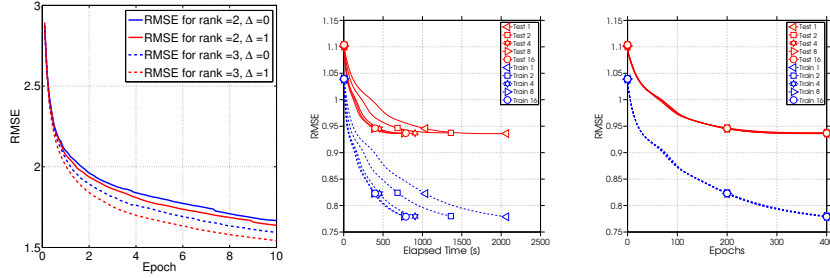


Figure 2: Left: The effect of adding inequalities ( $\Delta = 1$ ) to the equality-constrained problem ( $\Delta = 0$ ) on `smallnetflix`, for  $r = 2, 3$ ,  $\mu = 10^{-3}$ . Center and right: RMSE as a function of the number of iterations and wall-clock time, respectively, on a well-known  $480189 \times 17770$  matrix, for  $r = 20$  and  $\mu = 16$ .

uniformly at random, and the image was reconstructed using rank 50. The performance was measured in terms of PSNR, which is  $10 \log_{10}(255^2/E)$  for mean squared error  $E$ . Our approach with inequalities  $0 \leq Y_{i,j} \leq 255$  dominates all other approaches on 7 out of the 10 images. On the remaining 3 images, one would have to use the extrema of the observed elements, e.g., a subinterval of 12–246 for Barbara.

To illustrate the aggregate results further, we undertook the following experiment. We took a  $512 \times 512$  gray scale image (Lenna) and chose 50% of the pixels randomly, indexed as  $\mathcal{I}$ . Then, we ran Algorithm 1 for  $10^7$  serial iterations ( $|\hat{\mathcal{S}}| = 1$ ). We obtained solutions  $X_E(\text{rank})$  and  $X_{IN}(\text{rank})$ , where  $X_E(\text{rank})$  was obtained when we used only equality constrains ( $\mathcal{E} = \mathcal{I}, \mathcal{U} \equiv \mathcal{L} \equiv \emptyset$ ) and  $X_{IN}(\text{rank})$  was obtained when we used also inequality constrains ( $\mathcal{E} = \mathcal{I}, \mathcal{U} \equiv \mathcal{L} \equiv -\mathcal{I}, X^{\mathcal{L}} = \mathbf{0} \in \mathbb{R}^{512 \times 512}, X^{\mathcal{U}} = \mathbf{1} \in \mathbb{R}^{512 \times 512}$  and  $-\mathcal{I}$  is a set of all elements of  $X$  except those in  $\mathcal{I}$ ). Figure 3 shows for different  $\text{rank} \in \{30, 50, 100\}$  the best rank approximation obtained by SVD ( $X(\text{rank})$ ) and solutions  $X_E(\text{rank})$  and  $X_{IN}(\text{rank})$ . The benefit of obvious inequality constrains is nicely visible, e.g., at  $\text{rank} = 100$ , where the relative error of reconstruction is more than twice smaller. Further, the image is more smooth, upon visual inspection.

rank	$X(\text{rank})$	$X_E(\text{rank})$	$X_{IN}(\text{rank})$
30	 $\ X(\text{rank})\ _F = 223.9999$	 $Error = 13.1394$	 $Error = 12.6303$
50	 $\ X(\text{rank})\ _F = 224.6876$	 $Error = 18.2070$	 $Error = 13.1859$
100	 $\ X(\text{rank})\ _F = 225.2117$	 $Error = 39.1631$	 $Error = 15.2551$

Figure 3: Adding obvious constraints can help to get better solution. Error is defined as  $Error := \|X(rank) - X\|_F$ .

Table 1: Comparison with other solvers on the image recovery in terms of the peak signal-to-noise ratio (PSNR), citing the experiments of Wang et al. and adding results considering  $0 \leq Y_{i,j} \leq 255$  under “MACO”.

Instance / Algo.	SVT	SVP	SoftImpute	LMaFit	ADMIRA	JS	OR1MP	EOR1MP	MACO
Barbara	<b>26.9635</b>	25.2598	25.6073	25.9589	23.3528	23.5322	26.5314	26.4413	23.8015
Cameraman	25.6273	25.9444	26.7183	24.8956	26.7645	24.6238	27.8565	27.8283	<b>28.9670</b>
Clown	28.5644	19.0919	26.9788	27.2748	25.7019	25.2690	28.1963	28.2052	<b>29.0057</b>
Couple	23.1765	23.7974	26.1033	25.8252	25.6260	24.4100	27.0707	27.0310	<b>27.1824</b>
Crowd	<b>26.9644</b>	22.2959	25.4135	26.0662	24.0555	18.6562	26.0535	26.0510	26.1705
Girl	29.4688	27.5461	27.7180	27.4164	27.3640	26.1557	30.0878	30.0565	<b>30.4110</b>
Goldhill	28.3097	16.1256	27.1516	22.4485	26.5647	25.9706	28.5646	28.5101	<b>28.6265</b>
Lenna	28.1832	25.4586	26.7022	23.2003	26.2371	24.5056	28.0115	27.9643	<b>28.3581</b>
Man	<b>27.0223</b>	25.3246	25.7912	25.7417	24.5223	23.3060	26.5829	26.5049	26.5990
Peppers	25.7202	26.0223	26.8475	27.3663	25.8934	24.0979	28.0781	28.0723	<b>28.8469</b>

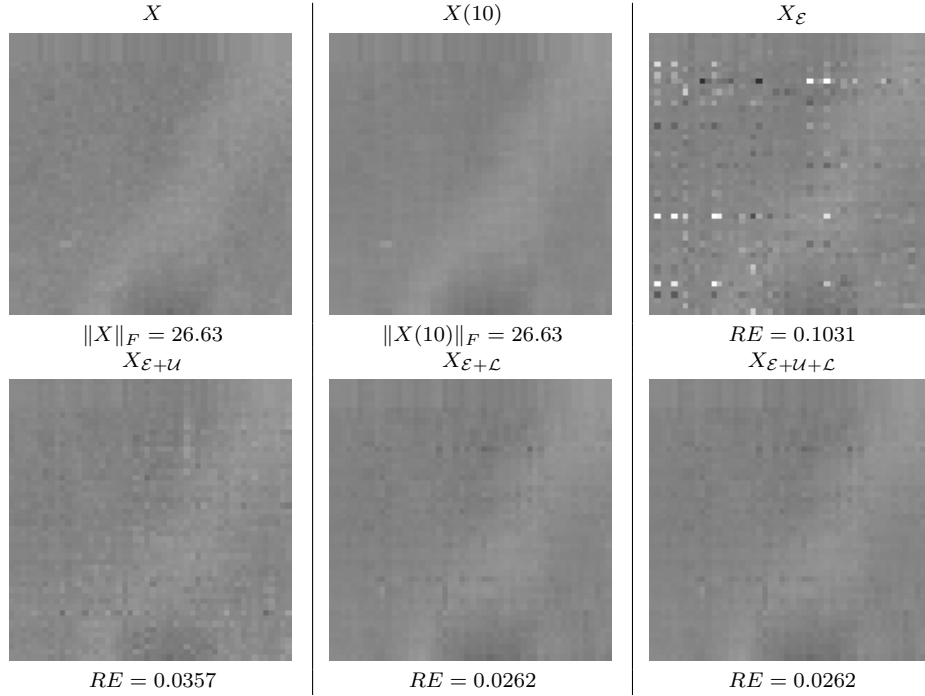


Figure 4: Original  $50 \times 50$  image, the best rank 10 approximation and reconstruction using Algorithm 1 with different settings. The  $RE$  is a relative error defined as  $RE(X.) = \|X. - X(10)\|_F / \|X(10)\|$ .

Further, we took a  $50 \times 50$  image and sampled randomly 50% of pixels. (The image is the top-left corner of the Lenna image.) Figure 4 shows the original image  $X$  and the best rank 10 approximation  $X(10)$ . The solutions  $X_E$ ,  $X_{E+U}$ ,  $X_{E+L}$  and  $X_{E+U+L}$  were obtained by running Algorithm 1 for  $3 \times 10^5$  serial iterations ( $|\hat{\mathcal{S}}| = 1$ ), where  $\mathcal{E}$  contains the observed pixels and  $\mathcal{U}$  and  $\mathcal{L}$  contains all other pixels. We have used  $X^{\mathcal{L}} = \mathbf{0}$  and  $X^{\mathcal{U}} = \mathbf{1}$ . The result again suggest that adding simple and obvious constrains leads to better low rank reconstruction and helps to keep reconstructed elements of matrix in expected bounds.

### 4.3 The Run-Time

Finally, in order to illustrate the run-time and efficiency of parallelization of Algorithm 1, Figure 2 (right) presents the evolution of RMSE over time on the well-known  $480,189 \times 17,770$  matrix of rank 20. There is an almost linear speed-up visible from 1 to 4 cores and marginally worse speed-up between 4 and 8 cores. Considering that most other algorithms proposed in the literature cannot cope with instances of this size, we cannot compare the performance

directly to SVT [6], SVP [16], SoftImpute [25], LMaFit [11], ADMiRA, [22], JS [15], OR1MP [45], EOR1MP [45], and similar. We can, however, compare the run-time on the  $512 \times 512$  instances, detailed in Table 1.

## 5 Conclusions

We have studied the matrix completion problem under interval uncertainty and an efficient algorithm, which converges to stationary points of the NP-Hard, non-convex optimisation problem, without ever trying to approximate the spectrum of the matrix. In our computational experiments, we have shown that even the seemingly most trivial inequality constraints are useful in a number of applications. This opens numerous avenues for further research:

- **Forecasting with Side Information:** A related application comes from the forecasting of seasonal data, e.g. sales. Let us assume that in process  $\{X_t\}$ , one knows  $k + 1 = \tau$  such that  $F_X(x_{t_1+\tau}, \dots, x_{t_k+\tau}) = F_X(x_{t_1}, \dots, x_{t_k})$  for the cumulative distribution function  $F_X(x_{t_1+\tau}, \dots, x_{t_k+\tau})$  of the joint distribution of  $\{X_t\}$  at times  $t_1 + \tau, \dots, t_k + \tau$ . One can then formulate the forecasting into the future as a matrix completion problem, where there the historical datum at time  $t$  is at row  $\lfloor t/\tau \rfloor$ , column  $t \bmod k$  specified by an equality or a pair of inequalities, and where inequalities represent side information. For example in sales forecasts, one often has bookings for many months in advance and knows that the sales for the respective months will not be less than the bookings taken.
- **Non-negative matrix factorization:** The coordinate descent algorithm for the problem (4) is easy to extend, e.g., toward non-negative factorization. It is sufficient to modify lines 7 and 13 in Algorithm 1 as follows:  $L_{i,\hat{r}} = \max\{0, L_{i,\hat{r}} + \delta_{i,\hat{r}}\}$ ,  $R_{\hat{r},j} = \max\{0, R_{\hat{r},j} + \delta_{\hat{r},j}\}$ . One could consider extensions beyond box constraints on the individual elements as well.
- **Auto-tuning  $\mu$ :** If we have some *a priori* bound on the largest eigenvalue of the matrix to reconstruct, let us denote it  $\zeta$ , then we can modify lines 7 and 13 in Algorithm 1 as follows  $L_{i,\hat{r}} = \max\{\min\{\zeta, L_{i,\hat{r}} + \delta_{i,\hat{r}}\}, -\zeta\}$ ,  $R_{\hat{r},j} = \max\{\min\{0, R_{\hat{r},j} + \delta_{\hat{r},j}\}, -\zeta\}$ .

We would be delighted to share our code with other researchers interested in these and related problems. Currently, the code is available from <http://optml.github.io/ac-dc/>. Should it become unavailable, for any reason, we encourage researchers to contact us.



## References

- [1] Carlos M. Alaíz, Francesco Dinuzzo, and Suvrit Sra. Correlation matrix nearness and completion under observation uncertainty. *IMA J. Numer. Anal.*, 35(1):325–340, 2013.
- [2] R.M. Bell and Y. Koren. Scalable collaborative filtering with jointly derived neighborhood interpolation weights. In *ICDM*, pages 43–52, Oct 2007.
- [3] Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. Optim.*, 20(4):1956–1982, 2010.
- [4] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *J. Assoc. Comput. Mach.*, 58(3):11, 2011.
- [5] Emmanuel J Candès and Yaniv Plan. Matrix completion with noise. *Proc. IEEE*, 98(6):925–936, 2010.
- [6] Emmanuel J Candès and Benjamin Recht. Exact matrix completion via convex optimization. *FoCM*, 9(6):717–772, 2009.
- [7] Yudong Chen, Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust matrix completion with corrupted columns. pages 873–880, 2011.
- [8] Yuxin Chen and Yuejie Chi. Spectral compressed sensing via structured matrix completion. In *ICML*, pages 414–422, 2013.
- [9] Alexander L Chistov and D Yu Grigor’ev. Complexity of quantifier elimination in the theory of algebraically closed fields. In *MFCS*, pages 17–31. Springer, 1984.
- [10] Maryam Fazel. *Efficient Algorithms for Collaborative Filtering*. PhD thesis, Stanford University, 2002.
- [11] Donald Goldfarb, Shiqian Ma, and Zaiwen Wen. Solving low-rank matrix completion problems efficiently. In *Allerton*, pages 1013–1020. IEEE, 2009.
- [12] J.P. Haldar and D. Hernando. Rank-constrained solutions to linear matrix equations using powerfactorization. *IEEE Signal Proc. Let.*, 16(7):584–587, July 2009.
- [13] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.
- [14] Nicholas JA Harvey, David R Karger, and Sergey Yekhanin. The complexity of matrix completion. In *SODA*, pages 1103–1111, 2006.
- [15] Martin Jaggi and Marek Sulovský. A simple algorithm for nuclear norm regularized problems. In *ICML*, pages 471–478, 2010.

- [16] Prateek Jain, Raghu Meka, and Inderjit S Dhillon. Guaranteed rank minimization via singular value projection. In *NIPS*, pages 937–945, 2010.
- [17] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *STOC*, pages 665–674, 2013.
- [18] R. Kannan, M. Ishteva, and Haesun Park. Bounded matrix low rank approximation. In *ICDM*, pages 319–328, Dec 2012.
- [19] Raghunandan H Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a few entries. *IEEE T. Inf. Theory*, 56(6):2980–2998, 2010.
- [20] Raghunandan Hulikal Keshavan. *Efficient Algorithms for Collaborative Filtering*. PhD thesis, Stanford University, 2012.
- [21] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, Aug 2009.
- [22] Kiryung Lee and Y. Bresler. Admira: Atomic decomposition for minimum rank approximation. *IEEE T. Inf. Theory*, 56(9):4402–4416, Sept 2010.
- [23] Guoyin Li, Alfred Ka Chun Ma, and Ting Kei Pong. Robust least square semidefinite programming with applications. *Comput. Optim. App.*, 58(2):347–379, 2014.
- [24] Jakub Mareček, Peter Richtárik, and Martin Takáč. Distributed block coordinate descent for minimizing partially separable functions. In *Numerical Analysis and Optimization*, pages 261–288, 2015. Springer Proceedings in Math. and Statistics, Vol. 134.
- [25] Rahul Mazumder, Trevor Hastie, and Robert Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *J. Mach. Learn. Res.*, 11:2287–2322, 2010.
- [26] Bhaskar Mehta, Thomas Hofmann, and Wolfgang Nejdl. Robust collaborative filtering. In *RecSys*, pages 49–56, New York, NY, USA, 2007. ACM.
- [27] Andriy Mnih and Ruslan Salakhutdinov. Probabilistic matrix factorization. In *NIPS*, pages 1257–1264, 2007.
- [28] Balas Kausik Natarajan. Sparse approximate solutions to linear systems. *SIAM J. Comput.*, 24(2):227–234, 1995.
- [29] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM J. Optim.*, 22(2):341–362, 2012.
- [30] Benjamin Recht, Christopher Ré, Stephen J. Wright, and Feng Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *NIPS*, pages 693–701, 2011.

- [31] Jasson D. M. Rennie and Nathan Srebro. Fast maximum margin matrix factorization for collaborative prediction. In *ICML*, ICML '05, pages 713–719, New York, NY, USA, 2005. ACM.
- [32] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Math. Program.*, 144(2):1–38, 2014.
- [33] Peter Richtárik and Martin Takáč. On optimal probabilities in stochastic coordinate descent methods. *Optim. Lett.*, pages 1–11, 2015.
- [34] Peter Richtárik and Martin Takáč. Distributed coordinate descent method for learning with big data. *J. Mach. Learn. Res.*, to appear, 2016.
- [35] Peter Richtárik and Martin Takáč. Parallel coordinate descent methods for big data optimization. *Math. Program.*, 156(1):433–484, 2016.
- [36] Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. Application of dimensionality reduction in recommender system—a case study. In *WebKDD*, 2000.
- [37] Shai Shalev-Shwartz, Alon Gonen, and Ohad Shamir. Large-scale convex minimization with a low-rank constraint. In *ICML*, pages 329–336, 2011.
- [38] Allen L. Soyster. Convex programming with set-inclusive constraints and applications to inexact linear programming. *Oper. Res.*, 21(5):pp. 1154–1157, 1973.
- [39] Nathan Srebro. *Learning with matrix factorizations*. PhD thesis, MIT, 2004.
- [40] Nathan Srebro, Jason Rennie, and Tommi S Jaakkola. Maximum-margin matrix factorization. In *NIPS*, pages 1329–1336, 2004.
- [41] Jared Tanner and Ke Wei. Normalized iterative hard thresholding for matrix completion. *SIAM J Sci. Comput.*, 35(5), 2013.
- [42] Rachael Tappenden, Peter Richtárik, and Burak Büke. Separable approximations and decomposition methods for the augmented lagrangian. *Optim. Methods Softw.*, 30(3):463–668, 2015.
- [43] Rachael Tappenden, Martin Takáč, and Peter Richtárik. On the complexity of parallel coordinate descent. *arXiv preprint arXiv:1503.03033*, 2015.
- [44] Ryota Tomioka, Taiji Suzuki, Masashi Sugiyama, and Hisashi Kashima. A fast augmented lagrangian algorithm for learning low-rank matrices. In *ICML*, pages 1087–1094, 2010.
- [45] Zheng Wang, Ming-Jun Lai, Zhaosong Lu, Wei Fan, Hasan Davulcu, and Jieping Ye. Rank-one matrix pursuit for matrix completion. In *ICML*, pages 91–99, 2014.

- [46] John Wright, Arvind Ganesh, Shankar Rao, Yigang Peng, and Yi Ma. Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. In Y. Bengio, D. Schuurmans, J.D. Lafferty, C.K.I. Williams, and A. Culotta, editors, *NIPS*, pages 2080–2088. Curran Associates, Inc., 2009.
- [47] Jieping Ye. Generalized low rank approximations of matrices. *Mach. Learn.*, 61(1-3):167–191, 2005.

**Acknowledgements** The authors are grateful for the reviews, which have helped them to improve both the presentation and contents of the paper. In addition, the first author acknowledges funding from the European Union Horizon 2020 Programme (Horizon2020/2014-2020), under grant agreement number 688380. The second author would like to acknowledge support from the EPSRC Grant EP/K02325X/1, *Accelerated Coordinate Descent Methods for Big Data Optimization*.